

AD-A135 201

THE DISTRIBUTED V KERNEL AND ITS PERFORMANCE FOR
DISKLESS WORKSTATIONS; STANFORD UNIV CA DEPT OF
COMPUTER SCIENCE D R CHERTON ET AL JUL 83
STAN CS 83 973 MDA903 80 C 0102

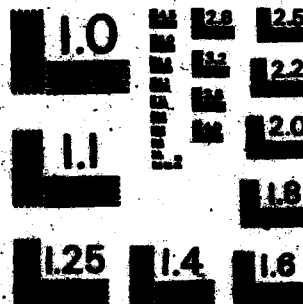
1/1

UNCLASSIFIED

1/6 17/2 NI



END
1/6



MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

July 1983

Report No. STAN-CS-83-973
Also numbered CSL 246

②

The Distributed V Kernel and Its Performance for Diskless Workstations

by

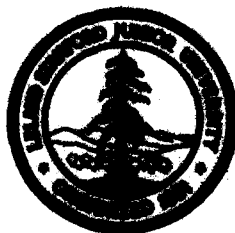
David R. Cheriton and Willy Zwanevool

Contract NSA-903-80-0-0103

Department of Computer Science

Stanford University
Stanford, CA 94305

APPROVED FOR PUBLIC RELEASE
DISTRIBUTION UNLIMITED



DTIC
ELECTE
DEC 1 1983
S A D

DTIC FILE COPY

88 : 10 04 015

AD-A185-201

The Distributed V Kernel and Its Performance for Diskless Workstations

David R. Cheriton and Willy Swenepoel

Computer Systems Laboratory
Departments of Computer Science and Electrical Engineering
Stanford University

The author

Abstract

The distributed V kernel is a microkernel-based system that provides efficient local and network transparent communication. It is primarily being used to do experiments of diskless workstations connected by a high-speed local network to a set of file servers. We describe a performance evaluation of the kernel, with particular emphasis on the cost of network file access. Our results show that over a local network:

1. Diskless workstations can access remote files with minimal performance penalty.
2. The V storage facility can be used to access remote files at comparable cost to a local disk-based storage facility.

It is concluded

that it is feasible to build a distributed system with all network communication using the V storage facility even when most of the network nodes have no permanent storage.

1. Introduction

The distributed V kernel is a microkernel-based system that provides efficient local and network transparent communication. The kernel consists of a set of processes that run on a set of nodes. It is designed to be used in a distributed system where nodes are connected by a high-speed local network. The kernel is designed to be used in a distributed system where nodes are connected by a high-speed local network. The kernel is designed to be used in a distributed system where nodes are connected by a high-speed local network.

Network interprocess communication is predominantly used for remote file access since most SUN workstations at Stanford are configured without a local disk.

This paper reports our experience with the implementation and use of the V kernel. Of particular interest are the experimental results of our approach, namely:

1. The use of diskless workstations with all secondary storage provided by remote file servers.
2. The use of a general purpose network transparent communication facility for support of special-purpose file access protocols and, in particular, the use of a Thrash-like transparent communication mechanism.

The most experimental approach is to configure workstations with a small local disk, using server-based file systems for primary storage. Diskless workstations, however, have a number of advantages including:

1. Lower hardware cost per workstation.
2. Greater responsiveness and simplicity of code with shared file servers.
3. Little or no memory or processing overhead on the workstations for the kernel and disk handling.
4. Tight coupling with application, reducing the distribution of data.

The major disadvantage is the overhead of processing all file access requests on the server. Our approach involves using the V storage facility to access remote files. This approach involves using the V storage facility to access remote files. This approach involves using the V storage facility to access remote files.

well as relatively easy to use, does not support application-level use of streaming.

These potential problems prompted a performance evaluation of our methods, with particular emphasis on the efficiency of file access. This emphasis on file access distinguishes our work from similar studies [10, 13]. The results of our study strongly support the idea of building a distributed system using distinct workstations connected by a high-speed local network to one or more file servers. Furthermore, we show that stream file access using the V kernel IPC facility is only slightly more expensive than a lower bound imposed by the basic cost of network communication. From this we conclude that relatively little improvement in performance can be achieved using protocols further specialized to file access.

2. V Kernel Interprocess Communication

The basic model provided by the V kernel is that of many small processes communicating by messages. A process is identified by a 32-bit globally unique process identifier or *pid*. Communication between processes is provided in the form of short fixed-length messages, each with an associated reply message, plus a data transfer operation for moving larger amounts of data between processes. In particular, all messages are a fixed 32 bytes in length.

The common communication scenario is as follows: A client process executes a *Send* to a server process which then completes execution of a *Receive* to receive the message and eventually executes a *Reply* to respond with a reply message back to the client. We refer to this sequence as a *message exchange*. The receiver may execute one of *MoveTo* or *MoveFrom* data transfer operations between the time the message is received and the time the reply message is sent.

The following sections describe the primitives relevant to this paper. The *kernel* refers to software in the V kernel named [1] for a complete description of the kernel facilities.

2.1. Primitives

Send(*message*, *pid*)

Send a message to the process identified by *pid*. The message is a 32-byte fixed-length message. The *message* parameter is a pointer to the message buffer. The *pid* parameter is a 32-bit process identifier. The *Send* operation is a blocking operation. The process calling *Send* will not return until the message has been received by the process identified by *pid*. The *Send* operation is a blocking operation. The process calling *Send* will not return until the message has been received by the process identified by *pid*.

Receive

(*pid*, *scope*) = *Receive*(*message*, *scope*, *reply*, *scope*)
Receive a message from the process identified by *pid*. The message is a 32-byte fixed-length message. The *message* parameter is a pointer to the message buffer. The *scope* parameter is a 32-bit process identifier. The *reply* parameter is a pointer to the reply message buffer. The *scope* parameter is a 32-bit process identifier. The *Receive* operation is a blocking operation. The process calling *Receive* will not return until the message has been received from the process identified by *pid*.

Reply(*message*, *pid*)

Send a 32-byte reply message to the process identified by *pid*. The message is a 32-byte fixed-length message. The *message* parameter is a pointer to the message buffer. The *pid* parameter is a 32-bit process identifier. The *Reply* operation is a blocking operation. The process calling *Reply* will not return until the message has been received by the process identified by *pid*.

ReplyWithSegment(*message*, *pid*, *segment*, *scope*, *scope*)

Send a reply message to the process identified by *pid* and transfer the data segment specified by *segment* and *scope* to the process identified by *scope*. The *message* parameter is a pointer to the message buffer. The *pid* parameter is a 32-bit process identifier. The *segment* parameter is a pointer to the data segment. The *scope* parameter is a 32-bit process identifier. The *ReplyWithSegment* operation is a blocking operation. The process calling *ReplyWithSegment* will not return until the message has been received by the process identified by *pid* and the data segment has been transferred to the process identified by *scope*.

MoveFrom(*segment*, *dest*, *src*, *scope*)

Copy data from the segment starting at *src* in the address space of *scope* to the segment starting at *dest* in the address space of *scope*. The *segment* parameter is a pointer to the data segment. The *dest* parameter is a pointer to the destination address. The *src* parameter is a pointer to the source address. The *scope* parameter is a 32-bit process identifier. The *MoveFrom* operation is a blocking operation. The process calling *MoveFrom* will not return until the data has been copied from the segment starting at *src* to the segment starting at *dest*.

MoveTo(*segment*, *dest*, *src*, *scope*)

Copy data from the segment starting at *src* in the address space of *scope* to the segment starting at *dest* in the address space of the process calling *MoveTo*. The *segment* parameter is a pointer to the data segment. The *dest* parameter is a pointer to the destination address. The *src* parameter is a pointer to the source address. The *scope* parameter is a 32-bit process identifier. The *MoveTo* operation is a blocking operation. The process calling *MoveTo* will not return until the data has been copied from the segment starting at *src* to the segment starting at *dest*.

GetFile(*segment*, *pid*, *scope*)

Return a pointer to the segment identified by *segment* in the address space of the process identified by *pid*. The *segment* parameter is a pointer to the data segment. The *pid* parameter is a 32-bit process identifier. The *scope* parameter is a 32-bit process identifier. The *GetFile* operation is a blocking operation. The process calling *GetFile* will not return until the segment has been located in the address space of the process identified by *pid*.

pid = *GetFile*(*segment*, *scope*)

Return the process identifier associated with *segment* in the address space of the process calling *pid*. The *segment* parameter is a pointer to the data segment. The *scope* parameter is a 32-bit process identifier. The *pid* parameter is a 32-bit process identifier. The *pid* operation is a blocking operation. The process calling *pid* will not return until the process identifier has been located in the address space of the process calling *pid*.

2.2. Operations

The V kernel provides a set of operations for moving data between processes. The *MoveFrom* and *MoveTo* operations are used to move data between segments in the address space of a process. The *Send* and *Reply* operations are used to send and receive messages between processes. The *GetFile* operation is used to get a pointer to a segment in the address space of a process. The *pid* operation is used to get the process identifier associated with a segment in the address space of a process.

The *MoveFrom* and *MoveTo* operations are used to move data between segments in the address space of a process. The *Send* and *Reply* operations are used to send and receive messages between processes. The *GetFile* operation is used to get a pointer to a segment in the address space of a process. The *pid* operation is used to get the process identifier associated with a segment in the address space of a process.

The *MoveFrom* and *MoveTo* operations are used to move data between segments in the address space of a process. The *Send* and *Reply* operations are used to send and receive messages between processes. The *GetFile* operation is used to get a pointer to a segment in the address space of a process. The *pid* operation is used to get the process identifier associated with a segment in the address space of a process.

must be allocated in the kernel and large amounts of data are transferred directly between user address spaces without user copies. Moreover, by virtue of the sparsity of the communication, the kernel's message buffers can be naturally allocated. As exemplified in Thrash, these factors make for a small, efficient kernel.

The V message primitives appear ill-suited in several ways for a network environment, at least on first observation. The short, fixed-length messages appear to make inefficient use of large packet sizes typically available on local networks. The synchronous nature of the message sending would seem to interfere with the true parallelism possible between separate workstations. And the economics of message buffering afforded by these restrictions in a single machine implementation are less evident in a distributed environment. Finally, the separate data transfer operations *MoveTo* and *MoveFrom* appear only to increase the number of remote data transfer operations that must be implemented in the distributed case.

However, our experience has been that the V message primitives are easily and efficiently implemented over a local network. Moreover, we have found that the messages of the primitives facilitated an efficient distributed implementation. The only major departure from Thrash was the explicit specification of segments in messages and the addition of the primitive *ReceiveWithSegment* and *ReplyWithSegment*. This extension was done for efficient page-level file access although it has proven useful under more general circumstances, eg. in pushing database string values to remote servers.

3. Implementation Issues

A foremost concern in the implementation of the kernel has been efficiency. Before describing some of the implementation details of the individual primitives, we first present aspects of the implementation that are common to the different operations of the kernel.

1. Remote operations are implemented directly in the kernel instead of through a process-level message layer. When the kernel receives a request from a remote process, it immediately copies the request into a kernel buffer. The kernel then sends the request to the appropriate process in the kernel. This is done by copying the request into a kernel buffer and then sending the request to the appropriate process in the kernel. This is done by copying the request into a kernel buffer and then sending the request to the appropriate process in the kernel.

2. The kernel uses a simple, efficient, and fast algorithm for sending and receiving messages. The kernel uses a simple, efficient, and fast algorithm for sending and receiving messages. The kernel uses a simple, efficient, and fast algorithm for sending and receiving messages. The kernel uses a simple, efficient, and fast algorithm for sending and receiving messages.

lowest functionality and local net performance, we have chosen not to include the dominant (but not) operation with any more overhead than is strictly necessary.

1. The synchronous request-response nature of a reply associated with each message is exploited to build reliable message transmission directly on an unreliable transport service, i.e. without using an extra layer (and extra packets) to implement reliable transport. The reply message serves as an acknowledgment as well as carrying the reply message data.
4. The mapping from process id to process location is aided by encoding a host specification in the process identifier. The kernel can then determine quickly whether a process is either local or remote, and in the latter case on which machine it resides.
5. There are no per-packet acknowledgments for large data transfers (as in *MoveTo* and *MoveFrom*). There is only a single acknowledgment when the transfer is complete.
6. File page-level transfer requires the optional number of packets (i.e. two) because of the ability to append short segments to messages using *ReceiveWithSegment* and *ReplyWithSegment*.

The following sections look at particular aspects of the implementation in greater detail.

3.1. Process Naming

V uses a global (but) naming space for specifying processes, in contrast to the local port naming used in DESOS [1] and Aegis [2]. Process identifiers are unique within the context of a host network. On the SUN workstation, it is natural for the V kernel to use 32-bit process identifiers. The high-order 16 bits of the process identifier serve as a logical host identifier which the low-order 16 bits are used as a locally unique identifier.

In the current 32-bit identifier implementation, the top 8 bits of the logical host identifier are the physical network address of the workstation, making the process identifier an external address mapping device. In the 16-bit implementation, a single byte high-order to network address. When there is no other way for the specified logical host, the message is dropped. Now "logical host/network address" correspondence can be determined from message received. However, with such small or difference of difference to use logical host identifier being hard problem.

The use of an explicit host field in the process identifier allows distributed processing of remote process identifier requests. The kernel can then determine quickly whether a process is either local or remote, and in the latter case on which machine it resides. The kernel can then determine quickly whether a process is either local or remote, and in the latter case on which machine it resides.

GenM uses network browsing to determine the mapping of a logical process identifier to a physical identifier if the mapping is not known to the local kernel. Any kernel handling the mapping are reported to the browsing engine. The addition of local and remote scrapes was required to discriminate between server processes that serve only a single workstation and those that serve the network.

When a process identifier is specified to `Send` with a logical host identifier different from that of the local machine, the local pid validation test fails and `Send` calls `NonLocalSend` which handles transmission of the message over the network.

THE UNIVERSITY OF CHICAGO PRESS

As in the local case, major connections arise with *AddressTo* and *AddressFrom* because, by their definition, there is no need for grouping or building of the data in the kernel. The V kernel sends the data directly from the source address space into the network interface, and directly from the network interface to the destination address space.

The message and data transfer primitives provide efficient communication of small amounts and large amounts of data, less than 32 bytes or more, out of network packets. However, perhaps the most serious problem is an information content of data that is not efficiently handled by the Thrift primitives when much is required to be sent or received.

[illegible]

Advantages of the Interview

1. The advantages of the Thoth IPC model with the enhanced performance of a WFE-type file access protocol.
2. Compatibility with efficient local operation. For business segments may be accessed directly if it is the most efficient space or if the recipient process operates a DMA device.
3. Use of Access/STASignature and Read/STASignature is authorized and recommended for maximum security. (See Table 1)

An expected objection to my solution is the apparent redundancy and self-redundancy of the savings provision. We find the redundancy may be oversteering when the persistent redundancy of communication and capitalization between client and power personnel. Also, it is not unexpected that there is some overlap between efficiency and slowness. Given the redundancy currently upon system services through such systems that provide a procedural impetus in the savings provision, it is not inappropriate to make their comparison in the practice level for efficiency.

We now turn to the discussion of the performance evaluation of the board. We first define the term *board quality* as a reasonably lower bound to the state of board composition. Subsequently we discuss the efficiency of the board, both in terms of economic policies and decisions.

[illegible][illegible]

1. The cost of private operations versus the cost of the corresponding public provision.
2. The cost of the extra costs of private versus government provision.

1. The first step in the process is to identify the problem. This involves gathering information about the situation and understanding the needs of the stakeholders involved.

2. Once the problem is identified, the next step is to develop a plan. This involves setting goals, identifying resources, and determining the steps that need to be taken to address the problem.

3. The third step is to implement the plan. This involves putting the plan into action and monitoring progress to ensure that the goals are being met.

4. Finally, the fourth step is to evaluate the results. This involves assessing the effectiveness of the plan and making adjustments as needed to improve the outcome.

network. The network gateway is a function of the processor, the network, the network interface and the number of bytes transferred. It is the critical time penalty for interpreting the network, because each software module that could otherwise transmit the data by passing pointers. The network penalty is obtained by summing the data to transmit a byte from the main memory of one computer to the main memory of another and vice versa and dividing the cost here for the equivalent by 2. The equivalent is required a large number of times for statistical purposes. The modules are implemented at the data link layer using the language used so that no program or program-switching overhead appears in the results. The assumptions of hardware transactions and low network utilization are good approximations of most real network environments.

However, primarily providing a more realistic minimum caloric intake than the transfer than that suggested by the physical account, would require it include the potential and expected interface steps. For instance, a 20 bit transfer can mean 1000 bits from one application to another in 100 milliseconds. However, the time to generate the packet in the interface at the sending end and the time to transfer the packet out of the interface at the receiving end are comparable to the time for transmission. Thus, the time for the transfer from the point of view of the communicating software modules is at least one or three times as long as the transfer by the physical transfer one.

Measurement of network quality was made using the experimental 2 dB network. In all measurements, the network was normally idle due to the considerable time at which measurements had to be made. Table 4-1 lists our measurements of the 2 dB network quality for the SMC configuration using the 2 and 40 dBm transmitters with those given in Table 4-2. The network measurements prove that the time for the data to be transmitted based on the physical bit rate of the medium, namely 2.54 Mbit/sec, is adequate.

CONFIDENTIAL

Year	Market Time	Market Supply	
		1970	1971
60	174	640	640
61	148	140	640
62	148	140	640
63	148	140	640
64	148	140	640
65	148	140	640

1. The first part of the document is a header section containing the following information:

a. Title: [Illegible text]
 b. Author: [Illegible text]
 c. Date: [Illegible text]

2. The second part of the document is a list of references or sources. The list contains several entries, each starting with a number followed by a period and a description of the source. The text is heavily obscured by noise and artifacts, making it difficult to read.

3. The third part of the document is a section of text that appears to be a summary or conclusion. It contains several paragraphs of text, which are also heavily obscured by noise and artifacts.

4. The fourth part of the document is a section of text that appears to be a list of footnotes or additional information. It contains several entries, each starting with a number followed by a period and a description of the source. The text is heavily obscured by noise and artifacts, making it difficult to read.

5. The fifth part of the document is a section of text that appears to be a list of footnotes or additional information. It contains several entries, each starting with a number followed by a period and a description of the source. The text is heavily obscured by noise and artifacts, making it difficult to read.

The difference between the network data, computed in the network data file, and the estimated network quality data is accounted for primarily by the parameter that is present and controls the quality and thus relative the quality of the other end. For instance, with a 100% type packet and an 80% type packet, the copy time from memory to the network interface and vice versa is roughly 1.50 milliseconds in each direction. This, of the total 6.95 milliseconds, 1.50 is copy time, 5.45 is network transmission time and 3 is (approximately) network congestion latency. If we transfer a 100% packet with double latency, the transmission time is less than 2 milliseconds while the copy times remain the same, reducing the parameter that 75 percent of the cost is the network latency. The importance of the parameter speed is also illustrated by the difference in network quality for the two processors measured in Table 4.1.

With our interface, the parameter is required to copy the packet into the interface for transmission and this of the interface as reception (with the interface providing bandwidth as buffer buffering). From the copy time given above, we might expect that a DMA interface would significantly improve performance. However, we would predict that a DMA interface would not result in better latency performance for the network. This is because the network interface is a copy of the packet from the network interface, allowing it to read from the packet data immediately in its own interface. With a DMA interface, the copy would be made into the network interface from DMA's own main memory. Similarly, as transmission, the network interface would copy the packet from the network interface. Thus DMA interface copies the packet to the processor in the computer area of memory, copying the packet from the network copy operation. Finally, there is no network interface in the knowledge of a network DMA interface. The network interface moves data faster than a 100% DMA interface. This parameter is used in the CPU contention. In general, the rate levels of a main network interface appear to be 10% affecting the main processor rather than spending up operation that will use of network operations.

5. Kernel Performance

Our first kernel performance study is the question of network interface contention. The kernel and performance is provided in some of the data in Table 4.1. The data shows that the network interface is the main factor in the network interface contention.

Measurement of network contention was done using a low-priority "background" process on each workstation that repeatedly copies a packet to or from the loop. All other processor activities follow the processor allocation to the packet. Thus, the packet that used per operation as a workstation is the elapsed time since the processor that allocated to the "background" process divided by 11, the number of operations executed.

Using 1000 calls per operation and time which gives or takes 10 milliseconds, our measurements should be accurate to about 10 milliseconds except for the effect of variation in network load.

5.1. Kernel Measurements

Table 5.1 shows the results of our measurements of average contention for the workstation with the kernel running on workstation using a 100% processor and measured by a 100% packet. This data shows that a kernel-based operation, which is just the kernel internal control of a kernel operation. The kernel-based operation and kernel give the operation time to work operation without delay and latency. The kernel-based operation shows the time difference between the kernel and network operation. The kernel-based operation gives the network delay to the network of the workstation in part of the kernel operation. The kernel-based operation shows the parameter time and for the operation on the two machines involved in the network operation of the operation. Table 5.2 gives the same measurements using a 100% packet. The data for both processors shows that the effect of the parameter speed has on network contention performance. As expected, the data for both operations, being significantly only on the processor speed, are 25 percent faster on the 25 percent faster processor. However, the data 25 percent improvement for some operation follows the parameter in the most significant performance, being in our operations and is an internal measurement by the network delay (at least on a highly loaded network).

A significant part of network contention is the packet transfer contention. This is the contention for the packet transfer. The data shows that the network interface is the main factor in the network interface contention. The data shows that the network interface is the main factor in the network interface contention. The data shows that the network interface is the main factor in the network interface contention.

Kernel Performance

Kernel Operation	Elapsed Time			Network Penalty	Processor Time	
	Local	Remote	Difference		Client	Server
GetTime	0.07	-	-	-	0.07	-
Send-Receive-Reply	1.00	1.18	2.18	1.00	1.79	2.30
MoveFrom: 1024 bytes	1.26	9.83	7.77	8.15	3.76	5.69
MoveTo: 1024 bytes	1.26	9.85	7.79	8.15	3.59	5.87

Table 5-1: Kernel Performance: 3 MB Ethernet and 80 MHz Processor (times in milliseconds)

Kernel Operation	Elapsed Time			Network Penalty	Processor Time	
	Local	Remote	Difference		Client	Server
GetTime	0.06	-	-	-	0.06	-
Send-Receive-Reply	0.77	2.54	1.77	1.30	1.44	1.79
MoveFrom: 1024 bytes	0.95	8.60	7.65	6.77	3.32	4.78
MoveTo: 1024 bytes	0.95	8.60	7.65	6.77	3.17	4.95

Table 5-2: Kernel Performance: 3 MB Ethernet and 10 MHz Processor (times in milliseconds)

view interprocess communication as transparent across machines when the speed ratio is so large. However, an alternative interpretation is to recognize that the remote operation adds a delay of less than 2 milliseconds, and that in many cases this delay is insignificant relative to the time necessary to generate a request in the server. Furthermore, the waiting at client workstation processor is tiny with the server load for only 1.04 milliseconds out of the total 2.14 milliseconds time taken by the 80 MHz processor). Then, one can afford the processor on one machine to, for instance, moving a server process to another machine if its request processing generally requires more than 0.07 milliseconds of processor time, i.e. the difference between the local time (Send-Receive-Reply time and the local processor time is 0.07 milliseconds).

5.4. Small-Program Traffic

The situation so far has been for a single pair of processes communicating over the network. Operating systems and network configurations should be able to support the network communication to applications with other problems. Some hardware is designed to handle multiple connections and some software is designed to handle the situation in a more sophisticated manner.

A pair of workstations connected by a network can be used to test the network communication. The network communication is tested by sending a message from one workstation to another. The message is sent over the network and the receiving workstation receives it. The message is then sent back to the sending workstation. This process is repeated several times to test the network communication.

network in this fashion. Unfortunately, our measurements of this scenario turned up a hardware bug in our 3 MB Ethernet interface, a bug which caused many collisions to go undetected and show up as corrupted packets. Thus, the response time for the 80 MHz processor workstation in this case is 3.4 milliseconds. The increase in time from 1.18 milliseconds is accounted for almost entirely by the detected and retransmitted data (usually one per 200 packets) from the hardware bug. With standard network interfaces, we believe that the network can support any reasonable level of message communication without significant performance degradation.

A more critical resource is processor time. This is especially true for workstations which are servers that need to be the focus of a number of network connections. The problem just faced on workstations that a workstation is loaded to a point where 500 message exchanges per second, independent of the number of clients. The number is substantially lower for the client workstation, particularly when a remote client for the server processing is involved. The table summarizes an example in Table 5-1.

The network communication is tested by sending a message from one workstation to another. The message is sent over the network and the receiving workstation receives it. The message is then sent back to the sending workstation. This process is repeated several times to test the network communication. The network communication is tested by sending a message from one workstation to another. The message is sent over the network and the receiving workstation receives it. The message is then sent back to the sending workstation. This process is repeated several times to test the network communication.

serving the client process we are measuring and otherwise idle. A later section considers multi-client load on the file server.

We first describe the performance of random page-level file access.

6.1. Page-level File Access

Table 6-1 list the times for reading or writing a 512 byte block between two processes both local and remote using the 10 NFS processor. The times do not include time to fetch the data from disk but do indicate expected performance when data is buffered in memory. A page read involves the sequence of kernel operations: *Send-Receive-ReplyWithSegment*. A page write is *Send-ReceiveWithSegment-Reply*.

Random Page-Level Access

Operation	Elapsed Time			Network Penalty	Processor Time	
	Local	Remote	Difference		Client	Server
page read	1.31	5.96	4.25	1.89	2.50	3.28
page write	1.31	5.69	4.29	1.89	2.50	3.32

Table 6-1: Page-Level File Access: 512 byte pages (times in milliseconds)

The columns are to be interpreted according to the explanation given for similarly labeled columns of Tables 3-1 and 3-2. Thus, the time to read or write a page using these primitives is approximately 1.5 milliseconds more than the disk access penalty for these operations.

There are several considerations that apply to the use of remote operations being faster than local operations. There are special cases of these described for simple message exchanges. First, the extra 4.2 milliseconds seen for remote operations is relatively small compared to the time cost of the file system operations itself. In particular, disk access time can be estimated at 20 milliseconds (assuming 1000000 bps) and the time to transfer data at 2.5 milliseconds. The total time to read or write a 512 byte block is 2.5 milliseconds and 4.2 milliseconds respectively, making the cost of the remote operations 20 percent higher than the local operations.

The second consideration is that the time to transfer data over the network is a significant factor. In this case, the time to transfer data over the network is 1.89 milliseconds. This is the time to transfer data over the network and is not the time to transfer data over the disk. The time to transfer data over the disk is 2.5 milliseconds. The time to transfer data over the network is 1.89 milliseconds. The time to transfer data over the disk is 2.5 milliseconds. The time to transfer data over the network is 1.89 milliseconds. The time to transfer data over the disk is 2.5 milliseconds.

The third consideration is that the time to transfer data over the network is a significant factor. In this case, the time to transfer data over the network is 1.89 milliseconds. This is the time to transfer data over the network and is not the time to transfer data over the disk. The time to transfer data over the disk is 2.5 milliseconds. The time to transfer data over the network is 1.89 milliseconds. The time to transfer data over the disk is 2.5 milliseconds.

The fourth consideration is that the time to transfer data over the network is a significant factor. In this case, the time to transfer data over the network is 1.89 milliseconds. This is the time to transfer data over the network and is not the time to transfer data over the disk. The time to transfer data over the disk is 2.5 milliseconds. The time to transfer data over the network is 1.89 milliseconds. The time to transfer data over the disk is 2.5 milliseconds.

The fifth consideration is that the time to transfer data over the network is a significant factor. In this case, the time to transfer data over the network is 1.89 milliseconds. This is the time to transfer data over the network and is not the time to transfer data over the disk. The time to transfer data over the disk is 2.5 milliseconds. The time to transfer data over the network is 1.89 milliseconds. The time to transfer data over the disk is 2.5 milliseconds.

difference between the client processor time for remote page access and for local page access, namely 1.3 milliseconds. A processor cost of more than 1.3 milliseconds per request can be expected from the estimation made earlier using LOCUS figures.

These measurements indicate the performance when file reading and writing use explicit segment specification in the message and *ReceiveWithSegment* and *ReplyWithSegment*. However, a file write can also be performed in a more basic Thrift-like way using the *Send-Receive-MoveFrom-Reply* sequence. For a 512 byte write, this costs 8.1 milliseconds; file reading is similar using *MoveTo*. Thus, the segment mechanism saves 3.5 milliseconds on every page read and write operation, justifying this extension to the message primitives.

6.2. Sequential File Access

Sequential file access is the predominant pattern for file activity in most systems. Efficient file systems exploit this behavior to reduce the effect of disk latency by prefetching file pages (read-ahead) and aggressively seeking modified pages (write-behind). File systems and file transfer protocols typically implement mechanisms to reduce the effect of network latency on sequential file access.

Using the V kernel remote communication between a workstation and a file server, the file server can implement read-ahead and write-behind to reduce the effect of disk latency. However, there is no streaming in the network IFC to deal with network latency.

Two factors suggest that streaming can be done without a local cache on the server. First, local disks have a low latency of approximately 20 milliseconds and a high speed of 1000000 bps. Second, network I/O is a highly performance sensitive operation. In the current implementation, further reducing the effect of network latency. The pattern of streaming data over the network is to send a buffer of data over the network and to receive a buffer of data over the network. The pattern of streaming data over the network is to send a buffer of data over the network and to receive a buffer of data over the network. The pattern of streaming data over the network is to send a buffer of data over the network and to receive a buffer of data over the network.

The third factor is that the time to transfer data over the network is a significant factor. In this case, the time to transfer data over the network is 1.89 milliseconds. This is the time to transfer data over the network and is not the time to transfer data over the disk. The time to transfer data over the disk is 2.5 milliseconds. The time to transfer data over the network is 1.89 milliseconds. The time to transfer data over the disk is 2.5 milliseconds.

many current time-sharing systems do, such program loading could achieve the same performance given in the table, independent of disk speed. Thus, we argue that *Advent* and *More/From* with large transfer units provide an efficient program loading mechanism that is as fast as can be achieved with the given hardware.

7. File Server Issues

File server performance is a critical issue for desktop workstations. Unfortunately, we do not yet have experience with a V kernel-based file server. Thus, this section describes what we believe are the key issues and estimates performance without providing conclusive data. In general, we view the processor as the key resource to consider in file server performance because, as argued earlier, the network bandwidth is plentiful and disk scheduling and buffering issues are identical to those encountered in conventional multi-user systems.

The number of workstations a file server can support can be estimated from processor requirements. If we estimate page read or write processing overhead as roughly 15 milliseconds for file system processing (from LOCUS) plus 15 milliseconds for kernel operation (from Table 6-1), a page request costs about 7 milliseconds of processor time. Program loading appears to cost about 300 milliseconds for an average 64 kilobyte program. Estimating that 50 percent of the file requests are page requests, the average request costs 26 milliseconds. Thus, a file server based on the SUN workstation processor could support about 35 file requests a second. From this we estimate that one file server can serve about 10 workstations adequately, but 20 or more active workstations would lead to excessive delays. However, a desktop workstation system can easily be extended to handle more workstations by adding more file server modules along the network until one reaches a bottleneck for less than 100 workstations.

For some programs, it is advantageous to view the file server processor requirements to estimate the performance of the server, rather than to find the processor has a minimum and subsequently add more page requests to the processor, causing for a slow start and doing a lot of file reads while waiting for the data requests. If they require very limited interaction with the user.

On this basis, a file server should have a general purpose processor capable of the ability to selectively execute some programs. The cost for the processor is \$1000 to \$1500. A further argument for using a general purpose processor is that the cost of a general purpose processor is less than the cost of a specialized processor. The cost of a specialized processor is \$1000 to \$1500. A further argument for using a general purpose processor is that the cost of a general purpose processor is less than the cost of a specialized processor. The cost of a specialized processor is \$1000 to \$1500. A further argument for using a general purpose processor is that the cost of a general purpose processor is less than the cost of a specialized processor. The cost of a specialized processor is \$1000 to \$1500.

program, i.e. it is transparent except for performance.

8. Measurements with the 10 Mb Ethernet

Our limited access to a 10 Mb Ethernet has precluded basing our measurements on this standard local network. However, some preliminary figures using the 10 Mb Ethernet indicate the effect of using a faster network and slightly faster network interfaces. First, the remote message exchange time is 2.71 milliseconds using an 8 MHz processor, roughly the time for the 10 MHz processor on the 3 Mb network and .5 milliseconds better than the 8 MHz processor on the 3 Mb network. Second, the page read time is 5.72 milliseconds. Finally, the program loading time is much improved, achieving 255 milliseconds for a 64 kilobyte load using 16 Kb transfer units. We have not identified to what degree the improvement is due to the faster network speed versus the difference in the network interface.

9. Related Work

There are a number of previous and concurrent efforts in providing communication mechanisms for distributed systems. For brevity, we compare our work with only a representative sample that characterizes the search for, and evaluation of, reliable models and implementations.

Spencer's remote reference study [13] considered the feasibility of implementing remote load and spare operations over a local network. Nelson's work on remote procedure calls [10] investigates network communication for procedure-based systems analogous to what the V kernel provides for message-based systems. Radich and Robinson's kernel [12] implements a message system with a number of features such as non-blocking message sending that are not provided by the V kernel. Finally, LOCUS [14] implements network communication into a UNIX-like system in the form of a network-based file system.

Our work has followed the general approach of Spencer's and Nelson's work in using a message-based form of communication in place of sending and receiving every program, kernel, or message performance. However, we share Spencer's concern that the system which requires globally shared memory is a distributed system, which is especially the case for the V kernel. The V kernel provides a strong degree of separation between processes and requires protected provision of services in a multi-user, multi-workstation environment by having independent communication to the kernel I/O processor.

Our approach differs from Nelson's primarily in our use of a message-based communication mechanism. This provides a more secure, reliable, and clear of the sharing differences in message. Furthermore, the V kernel provides a mechanism to have a shared processor, all messages by the addition of a shared processor and message object. Under such shared processor, the system can be extended to have a shared processor, all messages by the addition of a shared processor and message object. Under such shared processor, the system can be extended to have a shared processor, all messages by the addition of a shared processor and message object.

versus message-based systems, although it is not clear these differences result in any significant difference in overall performance.

The V kernel performance is roughly comparable to that of the software implementations developed by Spector and Nelson, allowing for the non-trivial differences in operation semantics and host processors. We would hypothesize that V kernel performance could be improved by a factor of 30 using microcode, similar to the improvement observed by Spector and Nelson for their primitives. Unfortunately, neither Spector nor Nelson provides results that afford a comparison with our file access results. In general, their work has concentrated on the speed of the basic mechanism and has not been extended to measure performance in a particular application setting.

In comparison to Accent, the V kernel provides a primitive form of message communication, and benefits accordingly in terms of speed, small code size and ability to run well on an inexpensive machine⁵ without disk or microcode support. For instance, Accent messages require an underlying transport protocol for reliable delivery because there is no client-level reply message associated with every *Send* as in the V kernel. We do not at this time have performance figures for Accent.

LOCUS does not attempt to provide applications with general network interprocess communication but exploits carefully honed problem-oriented protocols for efficient remote file access. It is difficult to compare the two systems from measurements available given the differences in network speeds, processor speeds and measurement techniques. However, from the specific comparisons with LOCUS presented earlier, we would expect overall file access performance for the V kernel to be comparable to LOCUS running on the same machines and network.

However, the memory requirements for the V kernel are about half that of LOCUS compiled for the PDP-11 and probably more like one fourth when LOCUS is compiled for a 32-bit processor like the 68000. Thus, for graphics workstations or process control applications, for instance, the V kernel would be more attractive because of its smaller size, real-time orientation and its provision of general interprocess communication. However, the V kernel does not provide all the functionality of the LOCUS kernel which includes that of the UNIX kernel and more. When required with V, these additional facilities must be provided by server processes executing either on client workstations or network server machines.

10. Conclusions

We conclude that it is feasible to build a distributed system using diskless workstations connected by a high-speed local network to one or more file servers using the V kernel IPC. In particular, the performance study shows that V kernel IPC provides satisfactory

performance despite its generality. Because the performance is so close to the lower bound given by the network penalty, there is relatively little room for improvement on the V IPC for the given hardware regardless of protocol and implementation used.

The efficiency of file access using the V IPC suggests that it can not only replace page-level file access protocols but also file transfer and remote terminal protocols, thereby reducing the number of protocols needed. We claim that V kernel IPC is adequate as a transport level for all our local network communication providing each machine runs the V kernel or at least handles the interkernel protocol. We do, however, see a place for these specific protocols in internetworking situations.

In addition to quantifying the elapsed time for various operations, our study points out the importance of considering processor requirements in the design of distributed systems. More experience and measurement of file server load and workstation file access behavior is required to decide whether file server processing is a significant problem in using diskless workstations.

The V kernel has been in use with the diskless SUN workstations, providing local and remote interprocess communication, since September 1982. It is currently 38 kilobytes including code, data and stack. The major use of the network interprocess communication is for accessing remote files. Our file servers are currently 6 VAX/UNIX systems running a kernel simulator and file server program which provides access to UNIX system services over the Ethernet using interkernel packets. A simple command interpreter program allows programs to be loaded and run on the workstations using these UNIX servers. Our experience with this software to date supports the conclusions of the performance study that we can indeed build our next generation of computing facilities [8] using diskless workstations and the V kernel.

Acknowledgements

We are indebted to all the members of the V research group at Stanford, which at this point includes two faculty members and roughly ten graduate students. In particular, we wish to thank Keith Lauts for his patient comments on a seemingly endless sequence of drafts and Tim Mann for his many contributions to the design and the implementation of the kernel. We would also like to thank the referees whose comments and suggestions helped to enhance the clarity of the paper.

References

1. F. Buxton, J.H. Howard and J.T. Manning. Task Communication in DEDACS. Proceedings of the 8th Symposium on Operating System Principles, ACM, November, 1977, pp. 23-31. Publication (Operating Systems Review 12(3)).
2. A. Susskind, F. Buxton, V. Paxon. The SUN Workstation Architecture. Tech. Rep. 82, Computer Science Laboratory, Stanford University, March, 1982.

⁵ In addition to the workstation used for this study, the V kernel runs on a PDP-11.

3. D.R. Cheriton, M.A. Malcolm, I.S. Melen and G.R. Sagar. "Thoth, a Portable Real-time Operating System." *Comm. ACM* 22, 2 (February 1979), 105-115.

4. D.R. Cheriton. Distributed I/O using an Object-based Protocol. Tech. Rept. 81-1, Computer Science, University of British Columbia, 1981.

5. D.R. Cheriton. *The Thoth System: Multi-process Structuring and Portability*. American Elsevier, 1982.

6. D.R. Cheriton, T.P. Mann and W. Zwankepoel. V-System: Kernel Manual. Computer Systems Laboratory, Stanford University.

7. Digital Equipment Corporation, Intel Corporation and Xerox Corporation. The Ethernet: A Local Area Network - Data Link Layer and Physical Layer Specifications, Version 1.0.

8. K.A. Lantz, D.R. Cheriton and W.I. Nowicki. Third Generation Graphics for Distributed Systems. Tech. Rept. STAN-CS-82-938, Department of Computer Science, Stanford University, February, 1983. To appear in *ACM Transactions on Graphics*.

9. R.M. Metcalfe and D.R. Boggs. "Ethernet: Distributed Packet Switching for Local Computer Networks." *Comm. ACM* 19, 7 (July 1976), 395-404. Also CSL-75-7, Xerox Palo Alto Research Center, reprinted in CSL-80-2.

10. B.J. Nelson. *Remote Procedure Call*. Ph.D. Th., Carnegie-Mellon U., 1981. published as CMU technical report CMU-CS-81-119.

11. G. Popok, B. Walker, J. Chow, D. Edwards, C. Kline, G. Rudisin, G. Thiel. LOCUS: A Network Transparent, High Reliability Distributed System. Proceedings of the 8th Symposium on Operating Systems Principles, ACM, December, 1981, pp. 169-177.

12. R. Rathid and G. Robertson. Accent: A Communication Oriented Network Operating System Kernel. Proceedings of the 8th Symposium on Operating Systems Principles, ACM, December, 1981, pp. 64-75.

13. A. Spitzer. "Performing Remote Operations Efficiently on a Local Computer Network." *Comm. ACM* 25, 4 (April 1982), 246-250.

14. D. Swinshart, G. McDowell and D. Boggs. WFS: A Simple Shared File System for a Distributed Environment. Proceedings of the 7th Symposium on Operating Systems Principles, ACM, December, 1979, pp. 9-23.



Accession For	
DTIC GRAFI	<input checked="" type="checkbox"/>
DTIC TAB	<input type="checkbox"/>
Unannounced	<input type="checkbox"/>
Justification	<input type="checkbox"/>
Distribution	
Availability Codes	
Avail and/or	Special

DATE
FILMED
8